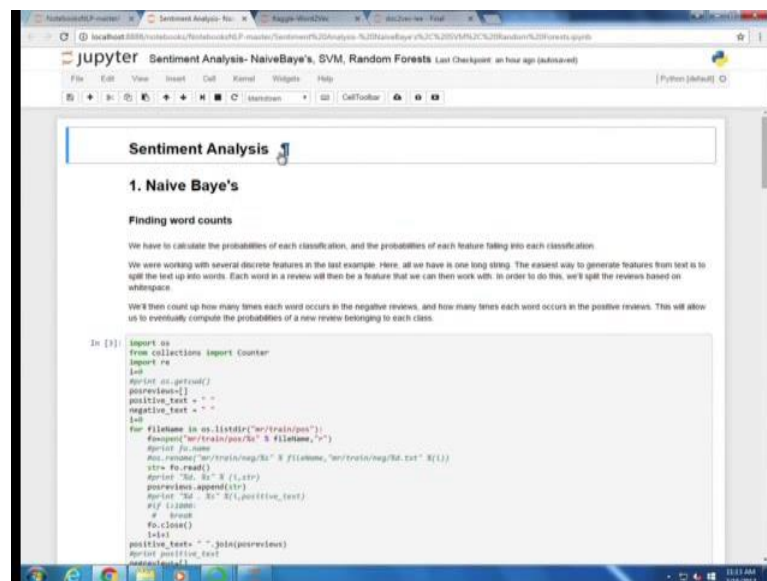**Natural Language Processing**
**Prof. Amrith Krishna**
**Department of Computer Science and Engineering**
**Indian Institute of Technology, Kharagpur**

**Lecture – 59**
**Tutorial**

Hello everyone, welcome to the tutorial section for the NLP course. So, today we will be looking into how to use or how to build a classifier for one of the most popular task in NLP sentiment analysis and we will be looking to how we create features and how we can build a supervised machine learning classifier for the purpose of sentiment analysis.
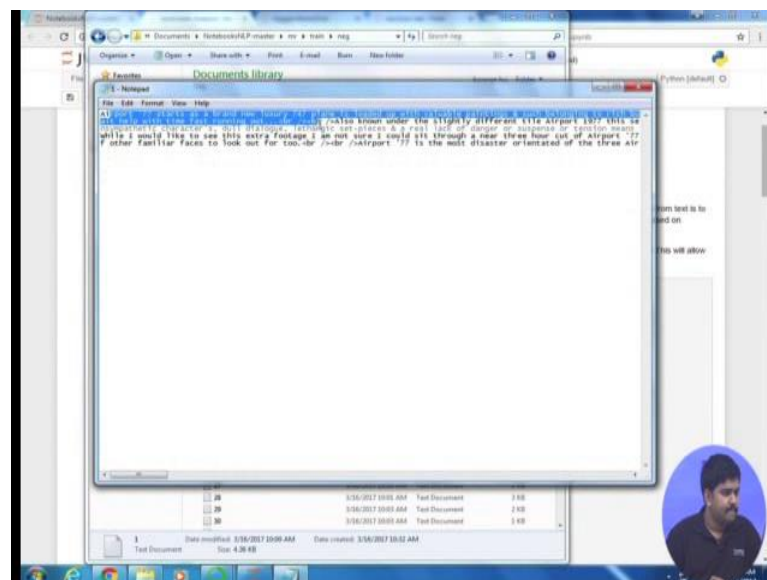
(Refer Slide Time: 00:24)



So, we will essentially look into how we can build a Naive Bayes classifier and after that we will just give you pointers towards how to build other classifiers. The notebook provided to you will have code snippets for building other type of classifiers, I will not be discussing it here.

So, sentiment analysis is essentially a task where you are given a document which can be a sentence or a set of sentences and you are given with and you have to label it whether it

belongs to a positive sentiment or a negative sentiment. This is the most simplistic representation or a task that we can think of sentiment analysis. Of course, as per the demands, this complexity of the task can be increased or like depending on the requirement from the customer or the user, there might be other levels of sentiment that might be applied here. So, here what we will be doing is that we will be taking about 5000 documents which are 5000 movie reviews and we see and we are already provided with the labels for those 5000 documents.

Here we can find about 2500 negative reviews and positive reviews about different films.

(Refer Slide Time: 02:05)



Now once this are provided, what we have to do is that we have to build a supervise classifier where we tag the positives as 1 and the negatives as 0 and then we build the classifier. So, here each document is your input and your models have to predict the sentiment of a new document when provided in the similar representation. So, we have to first convert each of those documents which is a training document in your vectorial representation. So, for that purpose, what we do here is that we take each positive and negative documents into separate data structure which is basically a list.
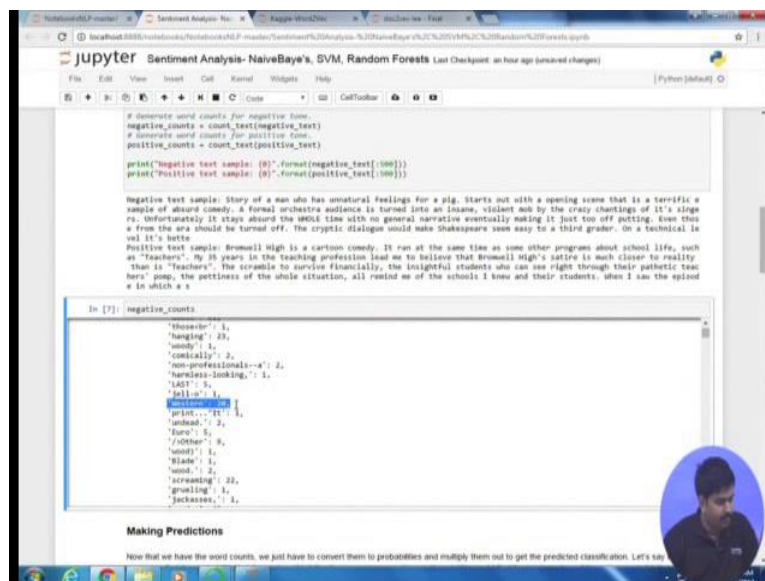
(Refer Slide Time: 03:28)



And once we do that what we do is that we count the number of individual words that is appearing in each of these sets.

We can find that since each document here is given as a separate file, we use the inbuilt OS library to traverse through each of the file given in the folder pos or the positive and then we append the content to a particular list. Once we do that we define a function called count text which basically looks for individual words in each of the document and also it keeps a track of the count or the number of time it is appearing.

(Refer Slide Time: 04:39)



We can just have a look of how the variable negative counts looks like, if you see here for negative the words pineapple, consider, daring, etcetera are appearing in different counts. So, we can find western appearing 20 times, screening appear 22 times, wooden appears 42 times and so on. So, here we have not done any pre processing for as it is so, far illustrative purposes. So, in this corpus may be the daring might be at separate word than this word because it has some extra symbols in it. These all need to be taken care when we go for a practical application. We do not do any case folding or any removal of other special symbols, but all these needs to be taken care.

So, once we have this separate dictionaries, now what we do is do here is that we just find out, how often a word is appearing in both positive and negative context. It is often possible that some of these words can appear in positive and negative context. So, the very basic idea here is that if a word appears more or less like equally in both the sets, they does not have much value, if a word words presents in one of these sets say positive context is queued with compared to what is happening with the negative set then that is more likely a positive sentiment word. If a word is more likely to belong to a negative set it is more likely a negative sentiment word, if it is more likely to belong to a positive set, it is a positive word.

(Refer Slide Time: 06:26)



We have the description here where we what we do is that we would find the total number of times a particular word that occur in the corpus and we just find the probability of it with its total count.

(Refer Slide Time: 06:50)



If we see what exactly is prob positive. So, what we do here is that we took to represent

the notion of cuteness or more likelihood of something being occurring in one set than the other, we convert them to a suitable probability distribution and what we find here is that the positive review count and negative review count. They are more or less equal they are like roughly 2500 a little bit here and there, but roughly they are equal. So, yes we can find what is the respective count.

We can find there is a slight difference of 300 also between both of them, but we essentially have more or less 2500 documents in each of those classes. Now so we can roughly say that they are like near to 0.5 in probability each, there is the number of documents that has prob positive sentiment and negative sentiment and are roughly half once we have that what remains is to find the individual probability of each word that is present. So, in order to do that, let us see, how this work, how this is calculated? So, print word and we will also find count, it seems to have a minor error in our work. So, indentation is something that is pretty important in python.

(Refer Slide Time: 11:41)



Here what we can find is that we can find the probability of each word belonging to a particular class. So, the function is defined in such a way that we first calculate all the prob, all the probability, all the words probability of belonging to the negative class then we find the probability of each of the word in the same document to belong to the

positive class. So, suppose our new sentence is good, nice, well and done, we will calculate first probability of those words belonging to the negative counts and for the negative probability set and then to the positive probability set.

If we see here, so done, good, well and nice, first when they are given to the negatives, we are getting a probability of 3.003 into 10 to the power of minus 6 as compared to 8.34 into 10 to the power minus 6, similarly for other values also like good negative has a probability of appearing with 4.5 into 10 to the power minus 5 where is to good in positive is higher than that. So, once we have this individual scores it take the product of this course as shown here and we find that the positive score for this document which convince this 4 words is 8.06 into the power 10 power minus 8 as compared to the negative score which is much lesser than this. So, we can say that overall sentiment for this review is positive.

Now let us feed for another sentence where the sentence is the movie was junk useless good for nothing sheer waste of time and money. So, it does not take much to think that what will be the label for this document, but what we have to see is that whether the system is able to capture that see if you find the sentiment of this particular document is negative and we can find that this particular document as a very huge difference in its positive score and negative score. We can find that the value here is in the order of 10 to the power minus 44 while the positive score is somewhere at in the value of 10 to the power of minus 47.

So, this means this particular document is labeled as negative and the system is able to capture both the cases of positive and negative documents are released for these 2 test document cases. So, this is how we internally calculate the values.
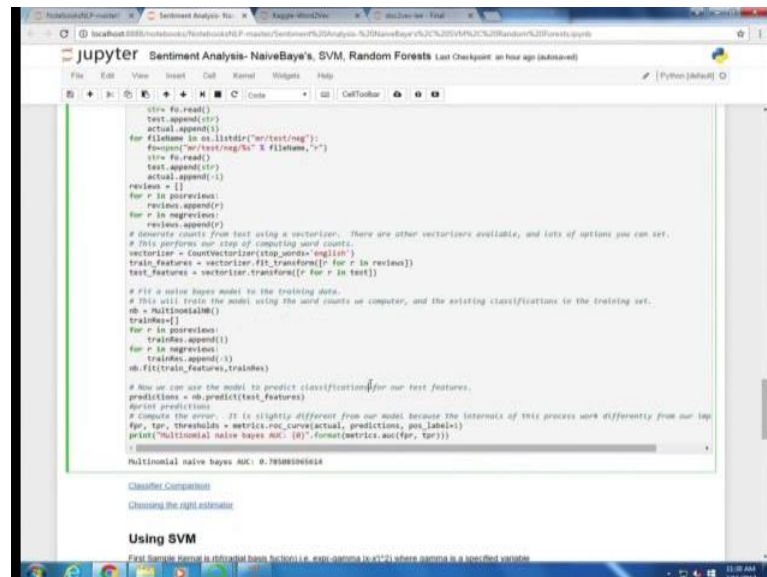
Now, how we can use a Naive Bayes; Naive Bayes classifier for a production ready purpose. So, here what we do is that we use the Scikit learn library they have a Naive Bayes classifier already implemented what we will do is that we will use this classifier to do the same task that is shown here, but we test it on a larger number of test data. So, while giving he input to the system we have to again convert it into a suitable vectorial representation that this particular library understands. So, we use the count vector is a function which basically represents a document in terms of the number of times a particular word is appearing in that document and it by default removes the stop words or it has other possible functions as well.

This particular function provides a whole range of facilities that makes our job much easier like we can often use different multi word patterns for multi word ranges like for example, if you have a word like hot dog though they are essentially separated by space, hot dog does not, when we take the individual words hot or dog, it does not capture the notion that is represented by hot dog. So, it is assumed to be a multi word expression. So, we will be; so these kind of Engram range functionality will help you to capture those multi word expressions and of course, we have to use some filtering criteria to remove the unnecessary ones and we convert both the training vectors and the test vectors into this vectorial representation. So, there is again a small difference in the function called

that we are doing fit transform and transform. So, fit and transform are actually two different methods and this Scikit learn function fit transform convert basically performs both the functions together. So, in fit we have to first build the vocabulary or the unique word.
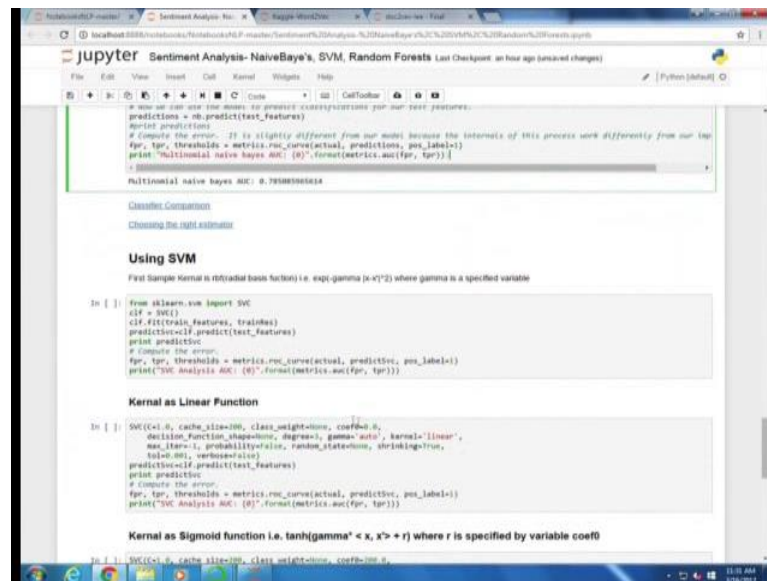
(Refer Slide Time: 16:48)



And then we convert each document to fit into any of to that vocabulary space in test features if a new word is appearing or that is called as out of vocabulary word that is ignored and only those counts are retained which are originally there in the vocabulary.

Once we have that we run this classifier and we also calculate the area under curve that gives us that basically give tells us about the false positive weight and the true positive weight, so for a binary classifier like this information helps us in identifying the performance of the classifier.

(Refer Slide Time: 18:14)



(Refer Slide Time: 18:25)



So, Scikit learn provides a host of classifiers it includes classifiers like SVM with multiple Kernal options or a simple classifiers like random forests. So, you can refer to any of the machine learning lectures to have an understanding of these classifiers for a quick check of what are the different classifiers that is provided.

(Refer Slide Time: 18:40)



(Refer Slide Time: 18:40)



You can have a look into the scikit learns official documentation. So, the link is already provided here. So, here the classifier essentially provides different classifiers and they show on one data set how each of this classifier is performing you can have a look into this in addition.
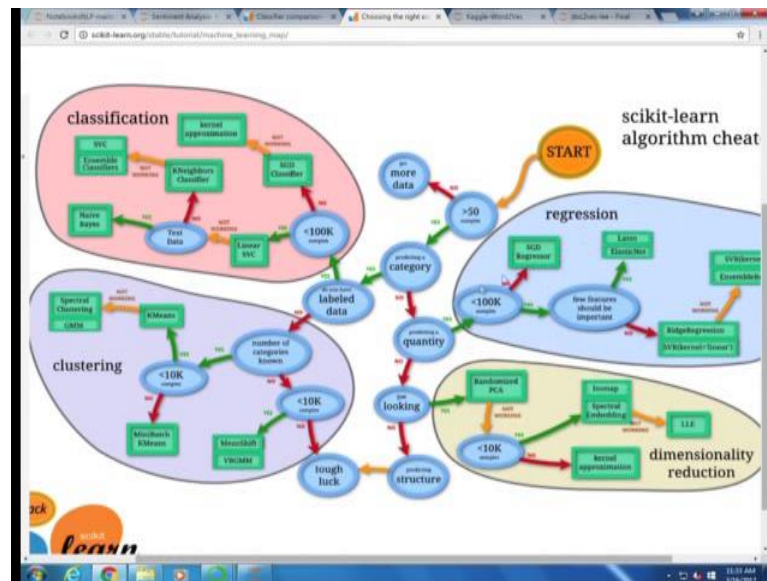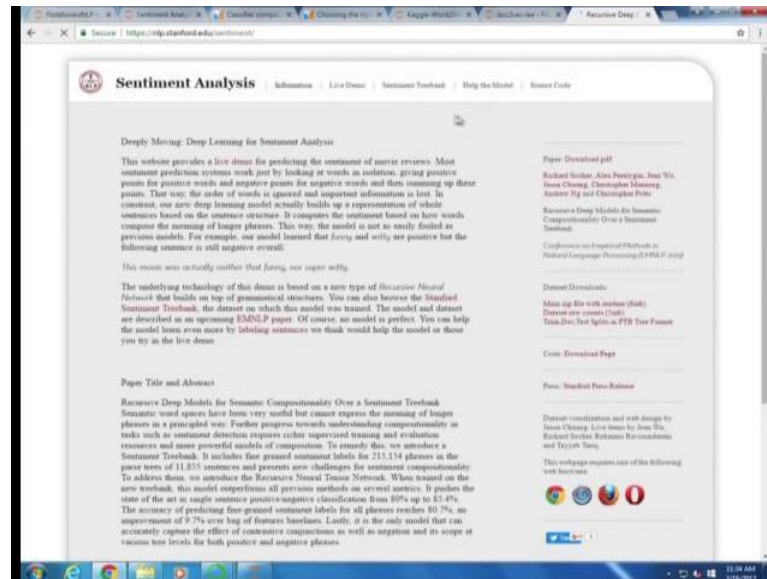
(Refer Slide Time: 19:09)



You can have a look into this quick scikit that will help you to find the relevant classifier or the estimator that you should be using. So, this scikit essentially provides all possible or it basically summarizes what are the different estimators that is available in scikit and you can see that when we do a supervise approach we should be either doing a classification or when we have where we have to classify or we have to predict a particular label.

In this case we were doing a classification where we predicting that whether something belongs to positive or negative. Suppose you wanted to predict the degree of positiveness or negativeness in a particular document you might have to model the same task as a regression task where it has a suite of library or suite of methods available here. In case of working with un supervised data like plastering or those kind of things that we have done in words disambiguation you might have to consider one of these clustering options or you can use the topic model that we have discussed in the previous tutorial.
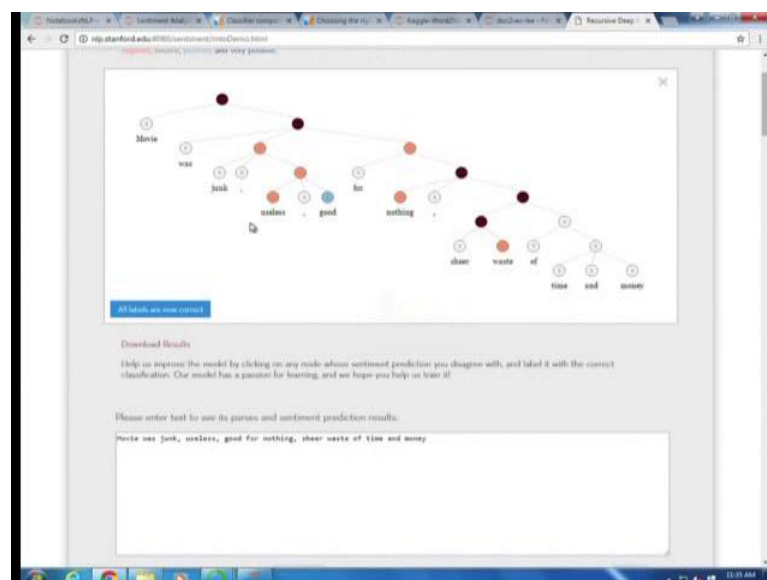
This essentially summarizes the class the supervised approaches for building a sentiment analysis tool.

(Refer Slide Time: 21:00)



You can also look into other standard sentiment analysis tools like the Stanford sentiment analyzer see if we give one of our sentences that we have used here for testing.

(Refer Slide Time: 21:37)



Here what we can find is that this again shows more or less negative sentiments to the classic the sentence. So, it basically shows each word to be either very negative or

negative and very few words to have a positive impact. So, that will be all for this tutorial, please try out other classifiers as well.

Thank you.